

## Lecture 11

Name: Tyler LaBonte

Professor: Santosh S. Vempala

# 1 Learning Halfspaces

Suppose our data domain is Euclidean  $\Omega = \mathbb{R}^d$  or boolean  $\Omega = \{0, 1\}^d$ . An important object is the halfspace, or linear threshold function. For  $\mathbf{w} \in \mathbb{R}^d$  and  $b \in \mathbb{R}$ ,

$$H(\mathbf{w}, b) = \{\mathbf{x} \in \Omega : \mathbf{w}^\top \mathbf{x} \geq b\}. \quad (1)$$

Many seemingly unrelated problem classes can be written as halfspaces, so they are pretty powerful. Let's see some examples. Suppose  $\mathbf{x} = (x_1, x_2, \dots, x_k)$  is an ordered set of booleans. Then a conjunction

$$x_1 \wedge \bar{x}_3 \wedge \dots \wedge x_k \quad (2)$$

with  $r$  negated terms may be written as the halfspace

$$\mathbf{w}^\top = (1, 0, 1, \dots, 1) \quad b = k - r. \quad (3)$$

Similarly, a disjunction

$$x_1 \vee \bar{x}_3 \vee \dots \vee x_k \quad (4)$$

with  $r$  negated terms may be written as the halfspace

$$\mathbf{w}^\top = (1, 0, 1, \dots, 1) \quad b = 1 - r. \quad (5)$$

One surprising equivalence is that of decision lists. The decision list

$$\text{if } x_1 = 1 \text{ then } 1 \text{ elif } x_2 = 0 \text{ then } 0 \text{ elif } \dots \text{ elif } x_k = 0 \text{ then } 0 \text{ else } 1 \quad (6)$$

may be written as the halfspace

$$2^k x_1 - 2^{k-1}(1 - x_2) + \dots + 2x_k + 1 \geq 1. \quad (7)$$

The exponentially decreasing weights ensure that the order of the decision list is preserved.

Note that, while there are infinitely many halfspaces in  $\mathbb{R}^d$ , they are still feasible to learn. To get an intuition for this, consider  $\Omega = \{0, 1\}^d$ . Say that halfspaces are distinct if they assign a different label to at least one element of the domain. Then, there are  $\binom{2^d}{d} \leq 2^{d^2}$  distinct halfspaces, which has a very reasonable logarithm of  $d^2$ . We will formalize this idea as VC-dimension in a

future lecture.

## 2 Perceptron Algorithm

Suppose we are in the mistake bound model and we want to learn the optimal halfspace  $\mathbf{w}^*$ . We can assume  $b = 0$ , since we can always append a 1 to  $\mathbf{w}^*$  and a  $-1$  to  $\mathbf{x}$ , thus working in  $\mathbb{R}^{d+1}$  with the same effect. We can also assume  $\|\mathbf{w}^*\|_2 = 1$  and  $\|\mathbf{x}\| \leq 1$ .

Define the label of a point in terms of the optimal halfspace:  $\ell(\mathbf{x}) = \text{sgn}(\langle \mathbf{w}^*, \mathbf{x} \rangle)$ . Then, we can learn a halfspace using the perceptron algorithm:

1. Initialize  $\mathbf{w} = \mathbf{0}$ .
2. On example  $\mathbf{x}$ , predict  $\text{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle)$ .
3. If mistake, update  $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{x}\ell(\mathbf{x})$ .

Intuitively, well-separated datasets should be more conducive to linear classifiers. This is quantified by the *margin*:  $\gamma = \min_{\mathbf{x}} |\langle \mathbf{w}^*, \mathbf{x} \rangle|$ .

### Theorem 2.1: Perceptron Mistake Bound

Suppose there exists  $\mathbf{w}^*$  with  $\gamma > 0$ . Then, the number of mistakes made by the perceptron algorithm is at most  $\frac{1}{\gamma^2}$ .

**Proof:** Consider the potential function

$$\cos(\mathbf{w}, \mathbf{w}^*) = \frac{\langle \mathbf{w}, \mathbf{w}^* \rangle}{\|\mathbf{w}\|}. \quad (8)$$

It starts out at zero. After one mistake,

$$\langle \mathbf{w}, \mathbf{w}^* \rangle \leftarrow \langle \mathbf{w} + \mathbf{x}\ell(\mathbf{x}), \mathbf{w}^* \rangle = \langle \mathbf{w}, \mathbf{w}^* \rangle + \ell(\mathbf{x})\langle \mathbf{w}^*, \mathbf{x} \rangle. \quad (9)$$

Since  $\ell(\mathbf{x})$  and  $\langle \mathbf{w}^*, \mathbf{x} \rangle$  have the same sign,  $\langle \mathbf{w}, \mathbf{w}^* \rangle$  increases by at least  $\gamma$  for each mistake. So, after  $t$  mistakes, it increases by at least  $\gamma t$ . For the denominator,

$$\langle \mathbf{w}, \mathbf{w} \rangle \leftarrow \langle \mathbf{w} + \mathbf{x}\ell(\mathbf{x}), \mathbf{w} + \mathbf{x}\ell(\mathbf{x}) \rangle \quad (10)$$

$$= \langle \mathbf{w}, \mathbf{w} \rangle + \langle \mathbf{x}, \mathbf{x} \rangle + 2\ell(\mathbf{x})\langle \mathbf{w}, \mathbf{x} \rangle. \quad (11)$$

The middle term is at most 1 by definition and the last term is negative because we made a mistake. Thus, it increases by at most 1 each mistake, and at most  $t$  after  $t$  mistakes. So,

$$\cos(\mathbf{w}, \mathbf{w}^*) \geq \frac{\gamma t}{\sqrt{t}} = \gamma\sqrt{t}. \quad (12)$$

If  $\gamma\sqrt{t} \leq 1$ , it must be that  $t \leq \frac{1}{\gamma^2}$ . ■

There is a generalization of this theorem with relaxed assumptions that obtains

$$t \leq \frac{\|\mathbf{w}^*\| \max \|\mathbf{x}\|}{\gamma^2}. \quad (13)$$

However,  $\gamma$  can be exponentially small even on the boolean cube, so it is not completely efficient. To solve this, there is a modified perceptron which makes  $\frac{\log n}{\sigma^2}$  mistakes and outputs a halfspace that answers correctly when the margin is above  $\sigma$ , and has no guarantees otherwise. Another way to solve this is by using linear programming.

### 3 The Kernel Trick

What if the threshold is not linear? Well, we can project to a higher-dimensional space where it is linear. For example, suppose the threshold is a quadratic,

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0 \implies \sum_{i,j} a_{ij} x_i x_j \geq 0. \quad (14)$$

Then, we let  $y_{ij} := x_i x_j$  to obtain the linear threshold  $\sum_{i,j} a_{ij} y_{ij} \geq 0$ . We now have  $\binom{d}{2}$  dimensions, but this is fine because perceptron gives a dimension-free guarantee!

#### Definition 3.1: Kernel

$K : \Omega \times \Omega \rightarrow \mathbb{R}$ ,  $K(\cdot, \cdot)$  is a legal kernel if  $\exists \phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$  such that  $K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$ . This is equivalent to the matrix  $\mathbf{K}$  being positive semidefinite.

Kernels help us compute inner products in projected high-dimensional spaces without actually knowing the mapping (which can be inefficient or have no closed form). Some kernels include  $\langle \mathbf{x}, \mathbf{y} \rangle^k$  and  $\exp(-\|\mathbf{x} - \mathbf{y}\|^2)$ .

To define the kernel perceptron, we first observe that in the perceptron algorithm, we can keep track of the set of mistakes instead of updating one at a time:  $\mathbf{w} = \sum_i \mathbf{x}_i \ell(\mathbf{x}_i)$  and prediction  $\langle \mathbf{w}, \mathbf{x} \rangle = \sum_i \langle \mathbf{x}_i, \mathbf{x} \rangle \ell(\mathbf{x}_i)$ . For kernels, we would have  $\mathbf{w} = \sum_i \phi(\mathbf{x}_i) \ell(\mathbf{x}_i)$ , but we don't want to actually compute  $\phi$ . Luckily, we can use the kernel:

$$\langle \mathbf{w}, \phi(\mathbf{x}) \rangle = \sum_i \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle \ell(\mathbf{x}_i) \quad (15)$$

$$= \sum_i K(\mathbf{x}_i, \mathbf{x}) \ell(\mathbf{x}_i). \quad (16)$$

Thus, we can learn a halfspace in  $\mathbb{R}^m$  without actually computing the mapping  $\phi$ .